

Benchmarking Large Language Models for News Summarization

Year: 2023 | Citations: 649 | Authors: *Tianyi Zhang, Faisal Ladhak, Esin Durmus*

Abstract

Large language models (LLMs) have shown promise for automatic summarization but the reasons behind their successes are poorly understood. By conducting a human evaluation on ten LLMs across different pretraining methods, prompts, and model scales, we make two important observations. First, we find instruction tuning, not model size, is the key to the LLM's zero-shot summarization capability. Second, existing studies have been limited by low-quality references, leading to underestimates of human performance and lower few-shot and finetuning performance. To better evaluate LLMs, we perform human evaluation over high-quality summaries we collect from freelance writers. Despite major stylistic differences such as the amount of paraphrasing, we find that LLM summaries are judged to be on par with human written summaries.